

# CSC449, HW#2, Kefu Zhu

## 1. Use of `dev` dataset

---

Since the I did not tune the learning rate in this assignment and I am using the SGD method to find the optimal weight vector, some updates during the iteration can potentially overshoot the local minima of the cost function.

Therefore, I treat the weight vector (model) from each iteration as separate models and created a variable `best_weights` to keep track the best weight vector that has been computed so far based on the test performance on `dev` dataset.

Rather than outputting the weight vector that results from the last iteration, I will return the best weight vector among all iterations when `--nodev` is not provided.

## 2. Model Training

---

As shown in the graph below, we can clearly see the model accuracy on both the `train` dataset and the `dev` dataset fluctuate a lot through the iterations, which also proves that the update is clearly overshooting a lot during the iteration

```
1  # Record the model performance (learning rate = 1)
2  performance_df_1 = model_performance(1,train_x,train_y,dev_x,dev_y)
3  # Plot the model performance (learning rate = 1)
4  fig, ax = plt.subplots(figsize = (15,7))
5  sns.lineplot(data = [performance_df_1.loc[:, 'train_accuracy'],
6                      performance_df_1.loc[:, 'dev_accuracy']])
7  ax.set_title('Model Performance over train and dev dataset (learning rate = 1)');
8  ax.set_xlabel('# of Iterations');
9  ax.set_ylabel('Accuracy');
10 ax.set_xticks(range(0,101,5));
11 ax.legend(['train', 'dev'])
```



After changing the learning rate from 1 to 0.01, the problem of overshooting is still significant. However, if we pay attention to the y-axis, we are able to reach higher accuracy than before.

```
1 # Record the model performance (learning rate = 0.01)
2 performance_df_001 = model_performance(0.01, train_x, train_y, dev_x, dev_y)
3 # Plot the model performance (learning rate = 0.01)
4 fig, ax = plt.subplots(figsize = (15,7))
5 sns.lineplot(data = [performance_df_001.loc[:, 'train_accuracy'],
6                     performance_df_001.loc[:, 'dev_accuracy']])
7 ax.set_title('Model Performance over train and dev dataset (learning rate = 0.01)')
8 ax.set_xlabel('# of Iterations');
9 ax.set_ylabel('Accuracy');
10 ax.set_xticks(range(0,101,5));
11 ax.legend(['train', 'dev'])
```



### 3. Appendix

**model\_performance** : Function used to produce the model performance graph

```
1 import pandas as pd
2 import seaborn as sns
3 import matplotlib.pyplot as plt
4
5 def model_performance(learning_rate, train_xs, train_ys, dev_xs, dev_ys):
6
7     # Initialize weight vector to be all zeros
8     weights = np.zeros(NUM_FEATURES)
9     # Initialize a dataframe to record model performance on train and dev dataset
10    df = pd.DataFrame(columns=['iter_num', 'train_accuracy', 'dev_accuracy'], dtype=
11
12    # Track the number of iterations
13    num_iter = 0
14    # Record the current number of iteration and model performance
15    df.loc[num_iter, 'iter_num'] = num_iter
16    df.loc[num_iter, 'train_accuracy'] = test_accuracy(weights, train_ys, train_xs)
17    df.loc[num_iter, 'dev_accuracy'] = test_accuracy(weights, dev_ys, dev_xs)
18    # While the number of iterations does not go beyond the maximum
19    while num_iter < 100:
20        # # Make a deep copy of weight vector before next iteration
```

```
21     # old_weights = np.array([w for w in weights])
22     # Loop through each pair of (x,y) in the training dataset
23     for x,y in zip(train_xs,train_ys):
24         # Classify current data point based on current weights
25         y_hat = np.sign(np.dot(np.transpose(weights),x))
26         # If we classify incorrectly
27         if y_hat!=y:
28             # Update the weight vector
29             weights = np.add(weights,np.array(learning_rate*y*x))
30
31     # Record the current number of iteration and model performance
32     df.loc[num_iter,'iter_num'] = num_iter
33     df.loc[num_iter,'train_accuracy'] = test_accuracy(weights, train_ys, train_xs)
34     df.loc[num_iter,'dev_accuracy'] = test_accuracy(weights, dev_ys, dev_xs)
35
36     # Increment the number of iteration
37     num_iter += 1
38     # Print the progress
39     if 100 <= 10: #
40         print('# of iterations: {}'.format(num_iter))
41     else:
42         if num_iter % round(100/10) == 0:
43             print('# of iterations: {}'.format(num_iter))
44
45     return df
```